

Sprachorientierte KI: Syntax und Parsing

- Syntax als Untersuchungsgegenstand
- Wortartendisambiguierung
- Phrasenstrukturgrammatiken
- Parsing mit Phrasenstrukturgrammatiken
- Restrikingierte Phrasenstrukturgrammatiken
- Unifikationsgrammatiken
- Constraint-basierte Grammatiken
- Robustes Parsing



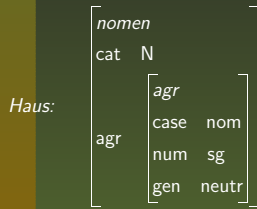
Constraint-basierte Grammatiken

- Getypte Merkmalstrukturen
- Verknüpfung von Merkmalstrukturen
- Systeme
- Constraints
- HPSG: Zeichen und Merkmale
- Prinzipien
- Dominanzschemata
- Lexikalische Regeln



Getypte Merkmalstrukturen

- Typen: jede Merkmalstruktur M erhält einen Typ t zugeordnet: M^t



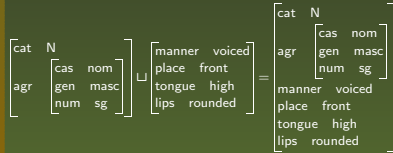
Getypte Merkmalstrukturen

- Erweiterung von Unifikation und Subsumtion auf getypte Merkmalstrukturen
- Subsumtion: $M_i^m \sqsubseteq M_j^n$ gdw. $M_i \sqsubseteq M_j$ und $m = n$
- Unifikation: $M_i^m \sqcup M_j^n = M_k^o$ gdw. $M_k = M_i \sqcup M_j$ und $m = n = o$



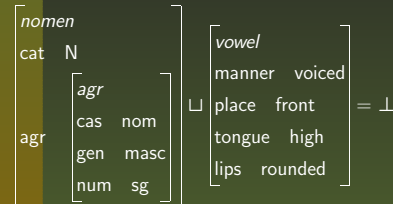
Getypte Merkmalstrukturen

- Effektivierung
 - Unifikation komplexer Merkmalstrukturen ist teuer
 - Vorentscheidung durch einfachen Typtest: Unifikation ist nur erforderlich, wenn Typtest erfolgreich war
- Programmierdisziplin
 - Beschränkung der seitlichen Erweiterbarkeit



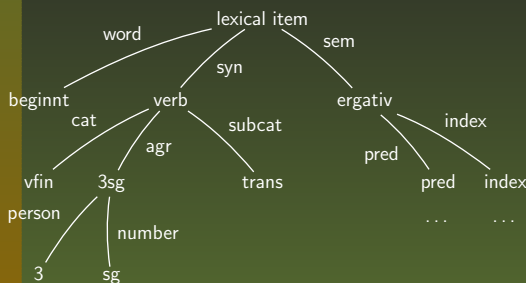
Getypte Merkmalstrukturen

- Ausschluss unbeabsichtigter Unifikationsresultate



Getypte Merkmalstrukturen

- grafische Interpretation: Typen als Knotenmarkierungen



Getypte Merkmalstrukturen

- Typenhierarchie:
 - partielle Ordnung über Typen: $\text{sub}(\text{verb}, \text{finit})$
 - hierarchische Abstraktion



Getypte Merkmalstrukturen

- Subsumtion für Typen:

$$m \sqsubseteq n \quad \text{gdw.} \quad \begin{cases} \text{sub}(m, n) \\ \text{sub}(m, x) \wedge \text{sub}(x, n) \end{cases}$$

- Unifikation für Typen:

$$m \sqcup n = o \quad \text{gdw.} \quad \begin{aligned} & m \sqsubseteq o \wedge n \sqsubseteq o \quad \text{und} \\ & \neg \exists x. m \sqsubseteq x \wedge n \sqsubseteq x \wedge x \sqsubseteq o \end{aligned}$$



Getypte Merkmalstrukturen

- Subsumtion für getypte Merkmalstrukturen:

$$M_i^m \sqsubseteq M_j^n \quad \text{gdw.} \quad \begin{aligned} & M_i \sqsubseteq M_j \quad \text{und} \\ & m \sqsubseteq n \end{aligned}$$

- Unifikation für getypte Merkmalstrukturen:

$$M_i^m \sqcup M_j^n = M_k^o \quad \text{gdw.} \quad \begin{aligned} & M_k = M_i \sqcup M_j \quad \text{und} \\ & o = m \sqcup n \end{aligned}$$



Getypte Merkmalstrukturen

- Zugehörigkeit (appropriateness)

- Spezifikation der zulässigen Attribute eines Typs, sowie der zulässigen Typen für deren Werte
- seitliche Erweiterbarkeit wird eingeschränkt
- Vererbung von Zugehörigkeitsfunktionen
- atomare Wertebelegungen als Typen: spezielle Typen definieren vollständig instanziierte Merkmalstrukturen
- Extremfall: alle Information in den Typen:
 - leere Merkmalstrukturen
 - Typhierarchie enthält alle vorausberechneten Unifikationsresultate



Getypte Merkmalstrukturen

- Typdefinition

- Vererbung von Merkmalstrukturen
- beliebig koreferente Strukturen
- entspricht den Templates, wenn Templatenname als Typbezeichnung verwendet wird



Verknüpfung von Merkmalstrukturen

- Verwendung der logischen Konnektoren

- Merkmalstrukturen als Beschreibungen für Denotatsmengen
 $\mathcal{I}(M) = \{x \mid M \text{ ist Beschreibung für } x\}$
- logische Konnektoren bilden Aussagen über die Zugehörigkeit zu Denotatsmengen
- aber zusätzlich immer Ermittlung der Beschreibung der resultierenden Denotatsmenge durch eine Merkmalstruktur



Verknüpfung von Merkmalstrukturen

- Konjunktion

- formal identisch mit der Unifikation
 $M_i \sqcup M_j = M_k$
 $\Leftrightarrow M_i \wedge M_j = M_k$
 $\Leftrightarrow \forall (x \in \mathcal{I}(M_k)). x \in \mathcal{I}(M_i) \wedge x \in \mathcal{I}(M_j)$
- Notation in CFS: $M_1 \sqcup M_2 \equiv \begin{bmatrix} M_1 \\ M_2 \end{bmatrix}$



Verknüpfung von Merkmalstrukturen

- Unifikation als Reduktion

- Zusammenfassen doppelter Attribute
- entspricht der Schnittmenge von Denotatsmengen
 $\mathcal{I}(M_i \sqcup M_j) = \{x \mid x \in \mathcal{I}(M_i) \cap \mathcal{I}(M_j)\}$



Verknüpfung von Merkmalstrukturen

- Disjunktion

- Generalisierung für disjunktive Merkmalstrukturen
- entspricht der Vereinigung von Denotatsmengen
- erfordert paarweise Unifikation der Disjunkte
- kombinatorisches Problem
- verzögerte Auswertung
- TMS-Verwaltung



Verknüpfung von Merkmalstrukturen

- Implikation:
 - Rückführung auf $M_i \rightarrow M_j \leftrightarrow \neg M_i \vee M_j$
 - erfordert die Komplementmengenbildung über Denotatsmengen $\mathcal{I}(M_i \rightarrow M_j) = D/\mathcal{I}(M_i) \cup \mathcal{I}(M_j)$
 - Denotatsmengen stehen nicht zur Verfügung, nur deren Beschreibungen durch Merkmalstrukturen
 - Approximation über das Pseudokomplement von M: allgemeinste Merkmalstruktur, deren Unifikation mit M scheitert
 - Pseudokomplement ist unter bestimmten Bedingungen nur eine Teilmenge des gewünschten Resultats
 - semantische Fundierung auf Feature-Logiken



Verknüpfung von Merkmalstrukturen

- approximative Definition der Implikation:

Sind M_i und M_j Merkmalstrukturen, so ist die Implikation $M_i \rightarrow M_j$ die allgemeinste Merkmalstruktur M_k , deren Unifikation mit M_i von M_j subsumiert wird:

$$M_i \rightarrow M_j = M_k \quad \text{gdw.} \quad M_j \sqsubseteq M_k \sqcup M_i.$$



Verknüpfung von Merkmalstrukturen

- Interpretation:

Wird eine Merkmalstruktur M von der linken Seite einer Implikation subsumiert, so kann die Information von der rechten Seite der Implikation zu M hinzuunifiziert werden

 - generelles Hinzufügen impliziter Information zu bestimmten Klassen von Merkmalstrukturen
 - Anwendung zur Formulierung von universalgrammatischen Prinzipien: Alle Merkmalstrukturen, die der linken Seite einer Implikation genügen, müssen auch mit der rechten Seite verträglich (unifizierbar) sein.



Verknüpfung von Merkmalstrukturen

- Ziel: kontextfreies Rückgrad beseitigen
 - keine individuellen Regeln
 - generelle Strukturierungsprinzipien
- Voraussetzung für die Suche nach dem Pseudokomplement: endliche Attributdomänen (Zugehörigkeitsfunktion)



Verknüpfung von Merkmalstrukturen

- Negation
 - Spezialfall der Implikation: $\neg M = M \rightarrow \perp$
 - Anwendung: Vermeiden langer Disjunktionen in der Lexikoninformation

Geld:
$$\left[\begin{array}{c} \text{cat} \quad \text{N} \\ \text{num} \quad \text{sg} \\ \text{agr} \quad \left[\begin{array}{c} \text{gen} \quad \text{neutr} \\ \{ [\text{cas} \quad \text{nom}] [\text{cas} \quad \text{dat}] [\text{cas} \quad \text{acc}] \} \end{array} \right] \end{array} \right] \quad \left[\begin{array}{c} \text{cat} \quad \text{N} \\ \text{num} \quad \text{sg} \\ \text{agr} \quad \left[\begin{array}{c} \text{gen} \quad \text{neutr} \\ \neg [\text{cas} \quad \text{gen}] \end{array} \right] \end{array} \right]$$

- Implementation als Makroexpansion

Voraussetzung: Zugehörigkeitsfunktion



Systeme

- Inferenzmaschinen, keine Parser
 - Ermittlung der speziellsten Merkmalstruktur für eine initiale partielle Beschreibung
 - Klassifikation
 - Typexpansion
 - Unifikation (Reduktion)



Systeme

- STUF (DÖRRE, EISELE, SEIFFERT, IBM Stuttgart/Heidelberg)
 - Stuttgarter Unifikationsformalismus
 - LILOG (STUF, STUF-II)
 - Verbmobil (STUF-III)
 - Typdefinition durch Zugehörigkeitsfunktion
 - keine Koreferenz in Typdefinitionen
 - Vorteil: Existenz der Denotatsmenge kann zur Übersetzungszeit überprüft werden
 - Nachteil: Spezifikation von Koreferenzen in externen Constraints



Systeme

- ALE (CARPENTER, CMU Pittsburgh)
 - attribute logic engine
 - Typdefinition durch Zugehörigkeitsfunktionen (intro) mit zusätzlichen Koreferenzconstraints (cons) und auswertbaren Prolog-Zielen (goal)
- TFS (EMELE, ZAJAC, Uni Stuttgart)
 - typed feature system



Systeme

- TDL (BACKOFEN, DFKI Saarbrücken)
 - type description language
 - Typdeklaration
 - gute Programmierumgebung
- ASL-Formalismus (EULER, Uni Hamburg)
 - nur Unifikator
- CFS (BÖTTCHER, KÖNYVES-TÓTH, GMD Darmstadt)
 - context feature system
 - Disjunktionsbehandlung über TMS



HPSG: Zeichen und Merkmale

- HPSG: head-driven phrase structure grammar (POLLARD, SAG 1987, 1994)
 - prinzipienbasierte Grammatik
 - universalgrammatischer Anspruch
- Zeichen
 - Zuordnung zwischen Zeichenkörper (Bezeichnendes, significant) und begrifflichem Konzept (Bezeichnetes, signifié)
 - Beschreibung durch Merkmalstrukturen



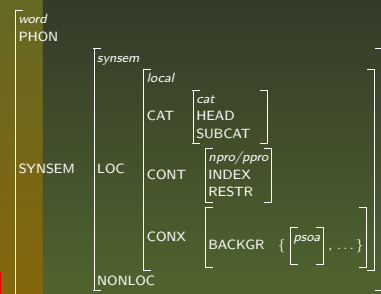
HPSG: Zeichen und Merkmale

- Klassifizierung der Zeichen in Typen
 - unterschiedliche Merkmale zur Beschreibung unterschiedlicher Zeichentypen
 - Merkmalstrukturen sind wohlgetypt: jeder Typ besitzt nur die zulässigen Attribute mit Werten vom geforderten Typ
 - Merkmalstrukturen sind total wohlgetypt: alle zulässigen Attribute sind auch spezifiziert
- atomare Merkmalstrukturen sind Typen, für die keine Merkmale zulässig sind



HPSG: Zeichen und Merkmale

- lexikalische Zeichen



HPSG: Zeichen und Merkmale

- PHONOLOGY (PHON)
 - phonetische Form
 - nicht systematisch behandelt
 - approximiert durch orthographische Repräsentation



HPSG: Zeichen und Merkmale

- SYNTAX-SEMANTICS (SYNSEM)
 - frühere Version Trennung
 - SEM: Bezeichnetes
 - SYN: strukturelle Vermittlung zwischen PHON und SEM
 - jetzt: Zusammenfassung
 - SYNSEM-Objekte als Grundlage der Subkategorisierung
 - enge Abhängigkeit: z.B. Modellierung von Kongruenz
 - Sprachen mit natürlichem Genus: semantische Ebene
 - Sprachen mit grammatischem Genus: syntaktische Ebene



HPSG: Zeichen und Merkmale

- LOCAL (SYNSEM|LOC)
 - gemeinsame Information von Spur und Filler bei ungebundenen Abhängigkeiten
- NONLOCAL (SYNSEM|NONLOC)
 - Grundlage für ungebundene Abhängigkeiten



HPSG: Zeichen und Merkmale

- CATEGORY (SYNSEM|LOC|CAT)
 - Kategoriale Information (SYNSEM|LOC|CAT|HEAD) Wortarteninformation
 - substantive (*subst*): Nomen, Verb, Adjektiv, Präposition
 - functional (*funct*): Determiner, Marker (z.B. Complementizer)
 - unterschiedliche Zugehörigkeitsfunktionen für unterschiedliche Typen, z.B.
 - noun*: CASE
 - verb*: VFORM, AUX, INV
 - ...



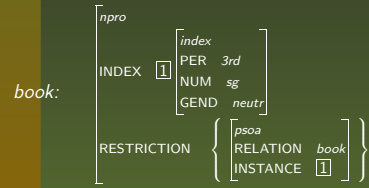
HPSG: Zeichen und Merkmale

- CATEGORY (SYNSEM|LOC|CAT)
 - Komplementforderungen (SYNSEM|LOC|CAT|SUBCAT)
 - Liste von *synsem*-Objekten
 - abzusättigende Valenzen einer lexikalischen Kategorie
 - einschließlich Subjekt (bzw. determinierende Subjekte einer NP)



HPSG: Zeichen und Merkmale

- CONTENT (SYNSEM|LOC|CONT)
 - Beitrag der lexikalischen Kategorie zur (kontext-unabhängigen, wörtlichen) Bedeutung einer Phrase
 - Referenzpotential
 - logische Form



HPSG: Zeichen und Merkmale

- *psoa*: parametrized state-of-affairs
 - Koreferenzen müssen für alle *psoa* verankert werden
 - Koreferenz mit einem Objekt vom Typ *book*
 - bei Bindung durch Allquantor (*every book*) Koreferenz mit allen Objekten, die die Restriktionen erfüllen



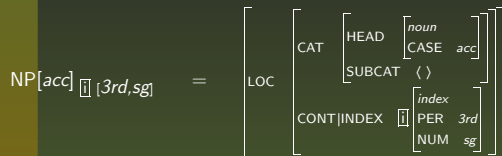
HPSG: Zeichen und Merkmale

- CONTEXT (SYNSEM|LOC|CONX)
 - kontextgebundener semantischer Beitrag der lexikalischen Kategorie
- BACKGROUND (SYNSEM|LOCAL|CONX|BACKGR)
 - Verankerungsbedingungen für Präsuppositionen und konventionale Implikaturen
 - natürliches Genus: Referent des englischen Pronomens *she* muß feminin sein



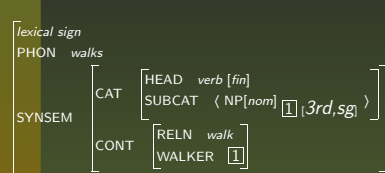
HPSG: Zeichen und Merkmale

- Notationsvereinfachung



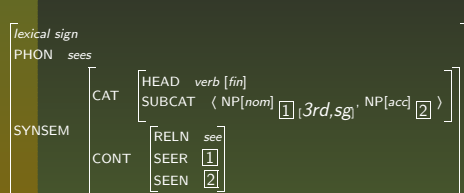
HPSG: Zeichen und Merkmale

- Beispiele für Lexikoneintragenen



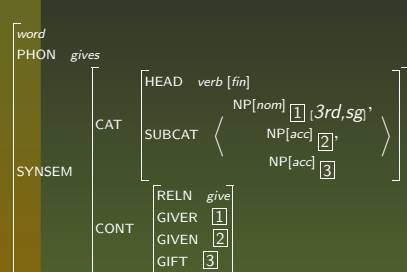
HPSG: Zeichen und Merkmale

- Beispiele für Lexikoneintragenen



HPSG: Zeichen und Merkmale

- Beispiele für Lexikoneintragenen



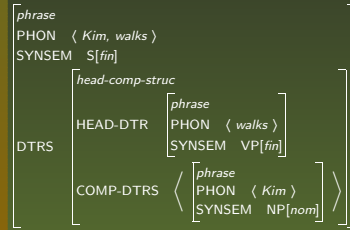
HPSG: Zeichen und Merkmale

- phrasale Zeichen
 - Zeichen des Typs *phrase*
 - zusätzliche Merkmale: Daughters, (Quantifier-Store)
 - wichtigster Spezialfall: *head-comp-struct* (Kopf-Komplement-Struktur)



HPSG: Zeichen und Merkmale

- DAUGHTERS (DTRS)
 - Konstituentenstruktur der Phrase
 - HEAD-DTR (*phrase*)
 - COMP-DTRS (Liste von Elementen des Typs *phrase*)



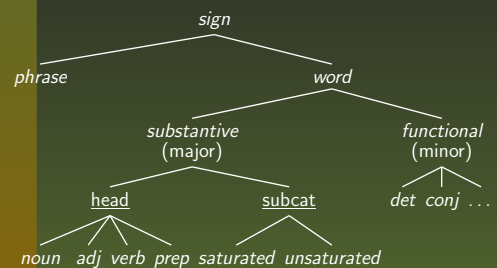
HPSG: Zeichen und Merkmale

- weitere Modellierungsgrundsätze
 - nichtderivationelle Beschreibung
 - structure sharing statt Bewegung
 - explizite Modellierung der Anordnung
 - keine Subjekt-Auxiliar-Inversion: verschiedene Anordnungsvarianten
 - keine Kopfbewegung von V^0 nach INFL



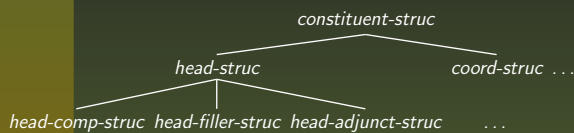
HPSG: Zeichen und Merkmale

- Zeichenhierarchie



HPSG: Zeichen und Merkmale

- Strukturtypen



Prinzipien

- Implikationen über getypten Merkmalstrukturen

$$\boxed{typ_1} \rightarrow \left[\begin{array}{l} X1 | \dots | XN \quad \boxed{1} \\ Y1 | \dots | YM \quad \boxed{1} \end{array} \right]$$

Hinzufügen typspezifischer Information

- Koreferenzen
- Wertebelegungen



Prinzipien

- universelle Prinzipien
 - Kopfmerkmalsprinzip
 - Subkategorisierungsprinzip
 - ...
- sprachspezifische Prinzipien
 - in der GB-Theorie: Parameter
 - aber: spekulativer Status der Parameter



Prinzipien

- Universalgrammatik
 - Linguistische Ontologie: universell verfügbare Typen, mit ihren Zugehörigkeitsfunktionen
 - Strukturschemata: begrenztes Inventar an universell verfügbaren Phrasentypen (schematische Dominanzregeln)
 - Kopf-Komplement-Strukturen, Kopf-Adjunkt-Strukturen, ...
 - Universelle Constraints
 - Kopfmerkmalsprinzip, Subkategorisierungsprinzip, ...



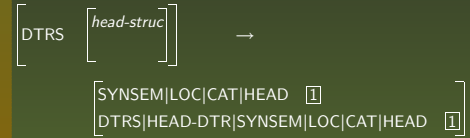
Prinzipien

- Sprachspezifische Grammatik
 - Lexikon (gegebenenfalls ergänzt durch die Anwendung lexikalischer Regeln)
 - Spezialisierung der linguistischen Ontologie
 - zusätzliche oder spezialisierte Strukturschemata



Prinzipien

- Kopfmerkmalsprinzip (Head-Feature-Principle)
 - Projektion der Kopfmerkmale an die Phrasenebene
 - Das HEAD-Merkmal einer Kopfstruktur ist mit dem HEAD-Merkmal seiner Kopftochter koindiziert.



Prinzipien

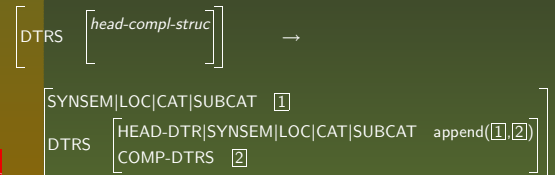
- Subkategorisierungsprinzip
 - Ordnung auf der SUBCAT-Liste: relative Obliquetheit
 - Subjekt ist nicht strukturell determiniert, sondern das Element der SUBCAT-Liste mit der geringsten Obliquetheit
 - Obliquetheithierarchie
 - Subjekt, primäres Objekt, sekundäres Objekt, oblique Präpositionalphrasen, Verbalkomplemente, ...
 - oblique Subkategorisierungsforderungen werden im Syntaxbaum zuerst abgebunden



Prinzipien

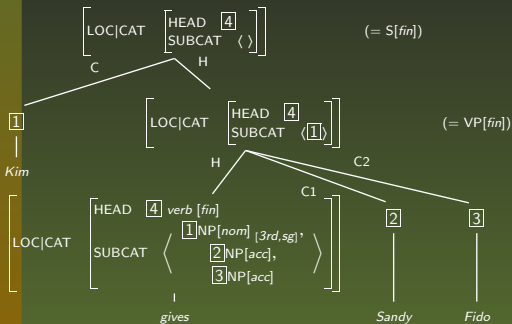
- Subkategorisierungsprinzip:

In einer Kopf-Komplement-Phrase ist der SUBCAT-Wert der Kopftochter gleich der Verkettung der SUBCAT-Liste der Phrase mit den SYNSEM-Werten der Komplementtöchter (geordnet nach steigender Obliquetheit).



Prinzipien

- Subkategorisierungsprinzip:



Prinzipien

- Semantikprinzip:

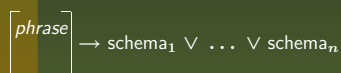
In einer Kopfphrase wird der CONTENT-Wert je nach Strukturtyp von der Adjunkttochter (*head-adjunct-struct*) oder der Kopftochter (sonst) auf die Phrase projiziert.
- Quantorenvererbung

Der Quantorenspeicher (QSTORE) einer Phrase ist die Vereinigung der Quantorenspeicher an den Tochterknoten, reduziert um die Quantoren, die am Mutterknoten selbst gebunden werden.
- SPEC-Prinzip



Dominanzschemata

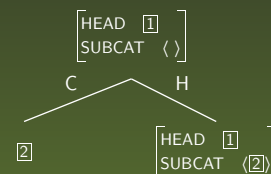
- X-Schemata: Bedingungen für wohlgeformte Teilbäume
 - $X^2 \rightarrow Y^2 \ X^1$ Spezifikatoradjunktion
 - $X^1 \rightarrow X \ Y^2$ Kopf-Komplementstruktur
- disjunkt spezifiziertes Prinzip der Universalgrammatik



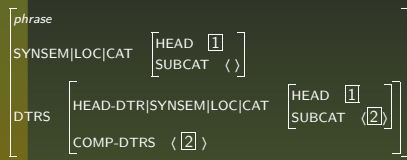
Dominanzschemata

- Schema 1:

eine gesättigte Phrase ($\left[\text{SUBCAT} \ \langle \rangle \right]$) mit einem DTRS-Wert vom Typ *head-comp-structure* in der der Wert des Merkmals HEAD-DTR ein phrasales Zeichen und der Wert von COMP-DTRS eine Liste der Länge eins ist.



Dominanzschemata



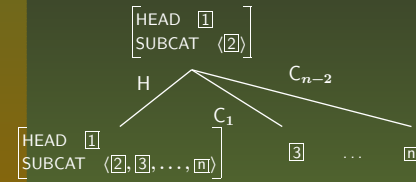
- direkte Konsequenz von Subkategorisierungs- und Kopfmerkmalsprinzip
- lizenzierte Phrasen
 - S → NP VP
 - NP → Det N¹



Dominanzschemata

■ Schema 2:

eine bis auf ein einziges Element abgesättigte Phrase mit einem DTRS-Wert vom Typ *head-comp-struct* in dem die Kopftochter ein lexikalisches Zeichen ist



Dominanzschemata

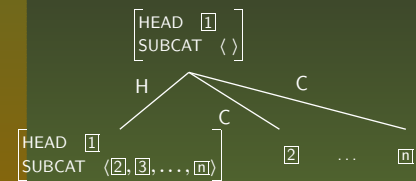
- lizenzierte Phrasen:
 - Verbalphrasen mit ihren Komplementen
- Komplexitätsebenen der \bar{X} -Theorie ersetzt durch Unterscheidung
 - lexikalisches / phrasales Zeichen bzw.
 - gesättigte / nichtgesättigte Phrase



Dominanzschemata

■ Schema 3:

eine gesättigte Phrase ([SUBCAT (< >)] mit einem DTRS-Wert vom Typ *head-comp-structure* und einem lexikalischen Kopf (Typ von HEAD-DTR ist *word*)



Dominanzschemata

- lizenzierte Phrasen:
 - "scrambling"-Strukturen: weitgehend freie Satzgliedanordnung (einschließlich Subjekt) z.B. Deutsch, Japanisch
- Schema 4: Head-Marker-Strukturen (*that John left*)
- Schema 5: Kopf-Adjunkt-Strukturen (z.B. Adjektive)
 - Grundprinzip: Adjunkte selektieren ihren Kopf

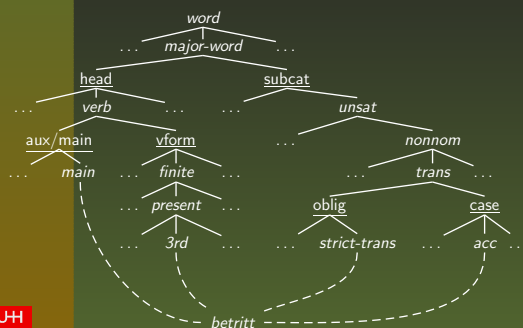


Lexikalische Regeln

- hierarchische Abstraktion in der lexikalischen Typhierarchie beseitigt "vertikale" Redundanz:
 - gemeinsame Information wird nur noch am gemeinsamen Supertyp gespeichert



Lexikalische Regeln



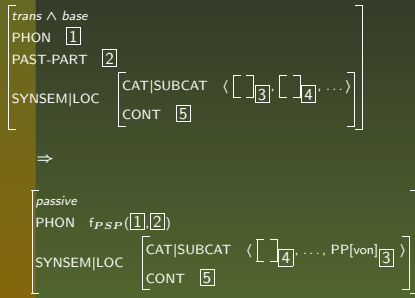
Lexikalische Regeln

- hierarchische Abstraktion beseitigt nicht die "horizontale" Redundanz im Lexikon:
 - abgeleitete Formen einer Grundform besitzen gemeinsame Informationen, die sich weitgehend regulär aus der Beschreibung der Grundform berechnen lassen
 - Flexion: Nomen, Verben, Adjektive, ...
 - Derivation: Nominalisierung, Partizipbildung, *un*-Negation, ...
 - Polyvalenz
 - Aktiv → Passiv
 - Dativobjekt → Präpositionalobjekt
 - ihm schreiben* → *an ihn schreiben*



Lexikalische Regeln

■ Beispiel: Passivierung



Lexikalische Regeln

- lexikalische Regeln sind nichtmonoton
- fügen sich schlecht in das Grundsche ma einer constraintbasierten Grammatik ein
- verschiedene Versuche zur Vermeidung
 - Morphological Principle (HENSCHTEL 1990)
 - gut geeignet für konkatenative Flexion
 - schlecht geeignet für syntaktische und semantische Veränderungen bei der Derivation
 - LEBETH 1994: monotone Beschreibung der Partizip-II-Bildung im Deutschen

